

Nouvelle représentation concise exacte des motifs corrélés rares : Application à la détection d'intrusions

Souad Bouasker, Tarek Hamrouni, Sadok Ben Yahia

Département des Sciences de l'Informatique, Faculté des Sciences de Tunis, Tunisie
{tarek.hamrouni, sadok.benyahia}@fst.rnu.tn

Résumé. La fouille des motifs corrélés qui sont très peu fréquents est une problématique de plus en plus intéressante dans la fouille de données. Dans ce cadre, les motifs corrélés rares selon la mesure de corrélation *bond* ont été étudiés dans un récent travail. La représentation concise exacte \mathcal{RMCR} de l'ensemble de ces motifs a été alors proposée. Toutefois, aucun algorithme n'a été proposé pour extraire cette représentation et aucune évaluation expérimentale de cette représentation n'a été réalisée. Dans ce papier ⁽¹⁾, nous introduisons l'algorithme RCPMINER d'extraction de \mathcal{RMCR} . Nous présentons également l'algorithme ESTMCR, d'interrogation de cette représentation ainsi que l'algorithme REGENERATIONMCR de dérivation de tous les motifs corrélés rares à partir de \mathcal{RMCR} . L'étude expérimentale réalisée montre des taux de compacité intéressants offerts par cette représentation. En outre, le processus de classification basé sur les règles génériques corrélées rares, dérivées à partir de \mathcal{RMCR} , a prouvé l'utilité de l'approche proposée dans le cadre de la détection d'intrusions.

1 Introduction et motivations

L'intégration des mesures de corrélation lors de l'extraction des motifs rares est une piste prometteuse en fouille de données. Elle permet, d'une part, d'améliorer la qualité des connaissances extraites en ayant un ensemble plus réduit contenant des motifs intéressants qui sont rares mais fortement corrélés. D'autre part, ceci renforce la qualité des règles d'association dérivées à partir de ces motifs corrélés rares. Par exemple, le motif composé par les items "Collier en or" et "Boucles d'oreilles" ou aussi celui composé de "Télévision" et "Lecteur DVD" correspondent à des motifs fortement corrélés mais peu fréquents dans les transactions d'une grande surface, et peuvent ainsi être omis dans un processus de fouille classique des motifs fréquents. L'utilité de tels motifs a été étudiée dans divers travaux tels que (Kim et al., 2011; Omiecinski, 2003; Segond et Borgelt, 2011; Surana et al., 2010; Xiong et al., 2006).

Dans la littérature, diverses approches d'extraction traitant de cette problématique ont été ainsi proposées. Nous citons, par exemple, l'approche décrite dans (Sandler et Thomo, 2010). Cette dernière est basée sur l'idée naïve d'extraire l'ensemble de tous les motifs fréquents pour

1. Ce travail propose une version étendue de l'article "Algorithmes d'extraction et d'interrogation d'une représentation concise exacte des motifs corrélés rares : Application à la détection d'intrusions", In Actes de la 12^{ième} Conférence Internationale Francophone Extraction et Gestion des Connaissances (EGC 2012), 31 Janvier - 03 Février 2012, Bordeaux, France.

un seuil minimal de support conjonctif, *minsupp*, très bas puis de filtrer ces motifs récupérés par la contrainte de corrélation. Cette opération est très coûteuse en temps de traitement et en consommation de la mémoire à cause de l'explosion du nombre de candidats à évaluer. Une autre stratégie d'extraction des motifs rares fortement corrélés, consiste à extraire l'ensemble de tous les motifs corrélés sans aucune intégration de la contrainte de support. Cette idée permet de récupérer les motifs corrélés qui sont très peu fréquents, cependant, elle est aussi coûteuse. Nous citons dans ce cadre les approches proposées dans (Ma et Hellerstein, 2001) et (Cohen et al., 2000). Il est important de noter que la contrainte monotone de rareté n'a été jamais incorporée dans la fouille afin de récupérer l'ensemble total des motifs rares fortement corrélés. En effet, les algorithmes proposés dans (Brin et al., 1997) et (Grahne et al., 2000), bien qu'ils permettent d'intégrer cette contrainte dans le processus de fouille, se limitent à l'extraction d'un sous-ensemble restreint composé uniquement des motifs minimaux valides *c.-à.-d.* satisfaisant l'ensemble de contraintes posées.

Dans (Bouasker et al., 2012), la représentation concise \mathcal{RMCR} des motifs corrélés rares associés à la mesure de corrélation *bond* (Omiecinski, 2003) a été proposée. D'un point de vue qualitatif, le choix de cette mesure a été effectué sur la base d'une étude détaillée de la littérature montrant son utilité dans le maintien de motifs intéressants (Ben Younes et al., 2012; Segond et Borgelt, 2011; Surana et al., 2010). D'un point de vue quantitatif, basée sur la notion clé de classe d'équivalence, cette représentation permet de ne présenter à l'utilisateur qu'un ensemble réduit de motifs tout en offrant la possibilité de dériver, si besoin, ceux non-retenus d'une manière simple et efficace. Toutefois, aucun algorithme n'a été proposé auparavant afin d'extraire une telle représentation. À cet égard, nous proposons, dans ce papier, un nouvel algorithme de fouille de la représentation \mathcal{RMCR} . Les algorithmes d'interrogation de cette représentation et de dérivation de l'ensemble total des motifs corrélés rares sont aussi présentés. En plus, nous décrivons les résultats obtenus prouvant les taux de compacité importants offerts par \mathcal{RMCR} ainsi que son apport dans la détection d'intrusions. Il est important de noter qu'aucun de ces algorithmes n'a été proposé et aucune expérimentation n'a été réalisée dans (Bouasker et al., 2012).

Le reste de ce papier est organisé comme suit : la section suivante présente l'ensemble des motifs corrélés rares et la représentation concise \mathcal{RMCR} qui lui est associée. Dans la section 3, nous introduisons l'algorithme RCPRMINER d'extraction de \mathcal{RMCR} . La section 4 est dédiée à la présentation de l'algorithme d'interrogation de \mathcal{RMCR} , tandis que la section 5 décrit le processus de régénération de l'ensemble de tous les motifs corrélés rares à partir de \mathcal{RMCR} . L'étude expérimentale est détaillée dans la section 6. L'application de la représentation \mathcal{RMCR} dans le cadre de la détection d'intrusions est illustrée dans la section 7. La conclusion et les perspectives de travaux futurs sont récapitulées dans la section 8.

2 Motifs corrélés rares : Définition et représentation concise

2.1 Notions de base

Nous commençons par définir d'abord une base de transactions.

Définition 1 (*Base de transactions*) Une base de transactions est représentée sous la forme d'un triplet $\mathcal{D} = (\mathcal{T}, \mathcal{I}, \mathcal{R})$ dans lequel \mathcal{T} et \mathcal{I} sont, respectivement, des ensembles finis de transactions (ou objets) et d'items (ou attributs), et $\mathcal{R} \subseteq \mathcal{T} \times \mathcal{I}$ est une relation binaire entre

	A	B	C	D	E
1	×		×	×	
2		×	×		×
3	×	×	×		×
4		×			×
5	×	×	×		×

TAB. 1 – Un exemple d'une base de transactions.

les transactions et les items. Un couple $(t, i) \in \mathcal{R}$ dénote le fait que la transaction $t \in \mathcal{T}$ contient l'item $i \in \mathcal{I}$.

Dans ce travail, nous nous sommes principalement intéressés aux itemsets comme classe de motifs. Nous distinguons trois types de supports correspondants à tout motif non vide X :

- **Le support conjonctif** : $SConj(X) = |\{t \in \mathcal{T} \mid \forall i \in X : (t, i) \in \mathcal{R}\}|$
- **Le support disjonctif** : $SDisj(X) = |\{t \in \mathcal{T} \mid \exists i \in X : (t, i) \in \mathcal{R}\}|$
- **Le support négatif** : $SNeg(X) = |\{t \in \mathcal{T} \mid \forall i \in X : (t, i) \notin \mathcal{R}\}|$

Exemple 1 Considérons la base de transactions illustrée par la table 1. Nous avons $SConj(AD) = |\{1\}| = 1$, $SDisj(AD) = |\{1, 3, 5\}| = 3$, et, $SNeg(AD) = |\{2, 4\}| = 2$ ⁽²⁾.

La fréquence conjonctive (*resp.* disjonctive et négative) est égale au support conjonctif (*resp.* disjonctif et négatif) divisé par $|\mathcal{T}|$. Dans la suite, nous allons utiliser les supports d'un motif. Comme nous nous intéressons aux motifs corrélés rares associés à la mesure de corrélation *bond* (Omiecinski, 2003), la définition suivante présente l'expression de *bond* telle que redéfinie dans (Ben Younes et al., 2012). Cette nouvelle expression permet de faire le lien entre la mesure *bond* et les supports conjonctif et disjonctif, cette mesure étant égale au rapport entre ces deux derniers.

Définition 2 (Mesure bond) Soit un motif non vide $X \subseteq \mathcal{I}$. La mesure *bond* de X est égale à :

$$bond(X) = \frac{SConj(X)}{SDisj(X)}$$

Ainsi, connaissant la valeur de *bond* et le support conjonctif d'un motif, il est aisé de dériver son support disjonctif et par conséquent son support négatif. Dans la sous-section suivante, nous présentons l'ensemble \mathcal{MCR} des motifs corrélés rares associés à la mesure *bond*.

2.2 L'ensemble \mathcal{MCR} des motifs corrélés rares

Les motifs corrélés rares ont été formellement définis dans (Bouasker et al., 2012) comme suit :

Définition 3 (Motifs corrélés rares) Étant donnés les seuils minimaux de support conjonctif et de corrélation *minsupp* et *minbond*, respectivement, l'ensemble \mathcal{MCR} des motifs corrélés rares est : $\mathcal{MCR} = \{X \subseteq \mathcal{I} \mid SConj(X) < minsupp \text{ et } bond(X) \geq minbond\}$.

Fouille d'une représentation concise des motifs corrélés rares

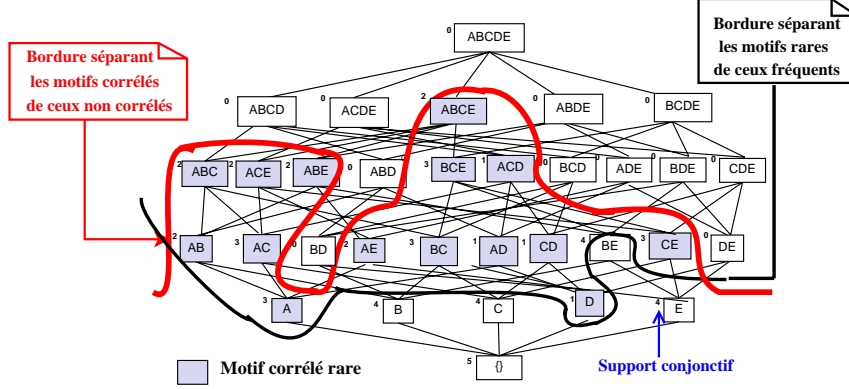


FIG. 1 – Espace des motifs corrélés rares pour $\text{minsupp} = 4$ et $\text{minbond} = 0,2$.

Exemple 2 Considérons la base illustrée par la table 1 pour $\text{minsupp} = 4$ et $\text{minbond} = 0,2$. L'ensemble \mathcal{MCR} est composé des motifs suivants où chaque triplet représente le motif, sa valeur de support conjonctif et sa valeur de bond : $\mathcal{MCR} = \{(A, 3, \frac{3}{3}), (D, 1, \frac{1}{1}), (AB, 2, \frac{2}{5}), (AC, 3, \frac{3}{4}), (AD, 1, \frac{1}{3}), (AE, 2, \frac{2}{5}), (BC, 3, \frac{3}{5}), (CD, 1, \frac{1}{4}), (CE, 3, \frac{3}{5}), (ABC, 2, \frac{2}{5}), (ABE, 2, \frac{2}{5}), (ACD, 1, \frac{1}{4}), (ACE, 2, \frac{2}{5}), (BCE, 3, \frac{3}{5}), (ABCE, 2, \frac{2}{5})\}$. Comme le montre cette figure, l'ensemble \mathcal{MCR} correspond aux motifs localisés en dessous de la bordure de la contrainte anti-monotone composée des motifs corrélés maximaux, et au dessus de la bordure de la contrainte monotone composée des motifs rares minimaux.

L'ensemble \mathcal{MCR} des motifs corrélés rares associés à la mesure bond résulte ainsi de la conjonction de deux contraintes de types opposés, à savoir la contrainte anti-monotone de la corrélation et la contrainte monotone de la rareté. Cette nature opposée des contraintes traitées rend complexe la localisation de l'ensemble des motifs corrélés rares. Ceci a motivé les auteurs dans (Bouasker et al., 2012) à introduire la représentation concise exacte \mathcal{RMCR} .

2.3 La représentation concise exacte \mathcal{RMCR}

La représentation concise exacte proposée constitue une réduction sans perte d'informations de l'ensemble \mathcal{MCR} . Pour cela, les auteurs ont recouru à la notion de bordure afin de délimiter l'espace associé à l'ensemble \mathcal{MCR} dans le treillis des motifs. Par ailleurs, l'ensemble des motifs corrélés rares a été ainsi partitionné en groupes disjoints, appelés "classes d'équivalence corrélées rares" en utilisant l'opérateur de fermeture f_{bond} (Ben Younes et al., 2012) associé à la mesure bond et défini comme suit.

Définition 4 (Opérateur f_{bond}) L'opérateur $f_{\text{bond}} : \mathcal{P}(\mathcal{I}) \rightarrow \mathcal{P}(\mathcal{I})$ associé à la mesure bond est défini comme suit : $f_{\text{bond}}(X) = X \cup \{i \in \mathcal{I} \setminus X \mid \text{bond}(X) = \text{bond}(X \cup \{i\})\}$.

2. Nous employons une forme sans séparateur pour les ensembles d'items : par exemple, AD représente l'ensemble $\{A, D\}$.

Chacune des classes d'équivalences induites par l'opérateur f_{bond} regroupe les motifs partageant les mêmes supports conjonctifs, disjonctifs et la même valeur de la mesure de corrélation $bond$. Les éléments maximaux des classes d'équivalence corrélées rares composent l'ensemble $\mathcal{MF}\mathcal{C}\mathcal{R}$ des motifs fermés corrélés rares et les éléments minimaux composent l'ensemble $\mathcal{MM}\mathcal{C}\mathcal{R}$ des motifs minimaux corrélés rares, qui ont été définis comme suit.

Définition 5 (Motifs fermés corrélés rares) L'ensemble $\mathcal{MF}\mathcal{C}\mathcal{R}$ des motifs fermés corrélés rares est défini par : $\mathcal{MF}\mathcal{C}\mathcal{R} = \{X \in \mathcal{M}\mathcal{C}\mathcal{R} \mid \forall X_1 \supset X : bond(X) > bond(X_1)\}$

Définition 6 (Motifs minimaux corrélés rares) L'ensemble $\mathcal{MM}\mathcal{C}\mathcal{R}$ des motifs minimaux corrélés rares est défini par : $\mathcal{MM}\mathcal{C}\mathcal{R} = \{X \in \mathcal{M}\mathcal{C}\mathcal{R} \mid \forall X_1 \subset X : bond(X) < bond(X_1)\}$.

Exemple 3 Soit la base illustrée par la table 1 pour $minsupp = 4$ et $minbond = 0,2$. Nous avons, par exemple, $f_{bond}(AB) = ABCE$, l'ensemble $\mathcal{MF}\mathcal{C}\mathcal{R} = \{A, D, AC, AD, ACD, BCE, ABCE\}$. Par ailleurs, l'ensemble $\mathcal{MM}\mathcal{C}\mathcal{R} = \{A, D, AB, AC, AD, AE, BC, CD, CE\}$. Il est intéressant de remarquer que les motifs A, D, AC et AD sont à la fois fermés et minimaux.

En se basant sur ces deux ensembles précédents, la représentation $\mathcal{RM}\mathcal{C}\mathcal{R}$ de l'ensemble $\mathcal{M}\mathcal{C}\mathcal{R}$ a été proposée.

Définition 7 (Représentation $\mathcal{RM}\mathcal{C}\mathcal{R}$) La représentation $\mathcal{RM}\mathcal{C}\mathcal{R}$ est définie comme suit : $\mathcal{RM}\mathcal{C}\mathcal{R} = \mathcal{MF}\mathcal{C}\mathcal{R} \cup \mathcal{MM}\mathcal{C}\mathcal{R}$.

Exemple 4 Considérons la base de transactions donnée dans la table 1, pour $minsupp = 4$ et $minbond = 0,2$. La représentation $\mathcal{RM}\mathcal{C}\mathcal{R} = \{(A, 3, \frac{3}{3}), (D, 1, \frac{1}{1}), (AB, 2, \frac{2}{5}), (AC, 3, \frac{3}{4}), (AD, 1, \frac{1}{3}), (AE, 2, \frac{2}{5}), (BC, 3, \frac{3}{5}), (CD, 1, \frac{1}{4}), (CE, 3, \frac{3}{5}), (ACD, 1, \frac{1}{4}), (BCE, 3, \frac{3}{5}), (ABCE, 2, \frac{2}{5})\}$.

Cette représentation a été prouvée dans (Bouasker et al., 2012) comme étant exacte, c.-à-d. permettant la régénération de tous les motifs corrélés rares sans perte d'informations. Par ailleurs, sa taille ne dépasse jamais celle de $\mathcal{M}\mathcal{C}\mathcal{R}$. En effet, $\mathcal{RM}\mathcal{C}\mathcal{R} = \mathcal{MF}\mathcal{C}\mathcal{R} \cup \mathcal{MM}\mathcal{C}\mathcal{R} \subseteq \mathcal{M}\mathcal{C}\mathcal{R}$.

Nous introduisons, dans ce qui suit, l'algorithme RCPRMINER permettant l'extraction de la représentation concise exacte $\mathcal{RM}\mathcal{C}\mathcal{R}$.

3 Algorithme RCPRMINER d'extraction de $\mathcal{RM}\mathcal{C}\mathcal{R}$

3.1 Description et pseudo code de l'algorithme RCPRMINER

L'algorithme RCPRMINER⁽³⁾, dont le pseudo-code est donné par l'algorithme 1, prend en entrée une base de transactions \mathcal{D} , un seuil minimal de support conjonctif $minsupp$ ainsi qu'un seuil minimal de corrélation $minbond$. Cet algorithme permet de déterminer, à partir du contexte \mathcal{D} , la représentation $\mathcal{RM}\mathcal{C}\mathcal{R}$ composée de l'ensemble $\mathcal{MM}\mathcal{C}\mathcal{R}$ des motifs minimaux corrélés rares et de l'ensemble $\mathcal{MF}\mathcal{C}\mathcal{R}$ des motifs fermés corrélés rares munis de leurs supports conjonctifs et de leurs valeurs de la mesure $bond$.

Le déroulement de l'algorithme RCPRMINER est illustrée par la figure 2. Cet algorithme se réalise en deux principales étapes. La première étape est dédiée à l'extraction, à partir de

3. Acronyme de Rare Correlated Patterns Representation Miner.

Fouille d'une représentation concise des motifs corrélés rares

\mathcal{D} , de l'ensemble \mathcal{MCM}_{\max} des motifs corrélés maximaux grâce à la procédure EXTRACTION_MCMAX (cf. ligne 4). Cette étape consiste à résoudre un problème classique permettant le repérage des éléments maximaux d'une théorie, les motifs maximaux associés à l'ensemble des motifs corrélés dans notre cas.

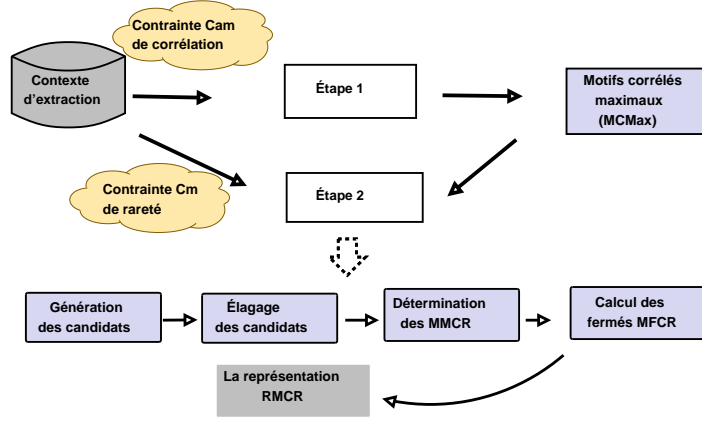


FIG. 2 – Schéma illustratif de déroulement de l'algorithme RCPRMINER.

La deuxième étape consiste à intégrer la contrainte de rareté ainsi que les motifs corrélés maximaux précédemment extraits dans un processus de fouille de la représentation \mathcal{RMCR} . À chaque itération de cette deuxième étape, l'ensemble \mathcal{CandP}_n est composé des candidats potentiels de taille n générés, moyennant la procédure APRIORI_GEN (cf. ligne 15), à partir des candidats retenus de taille $(n - 1)$. Les candidats de l'ensemble \mathcal{CandP}_n seront ainsi élagués (cf. ligne 11) selon différentes stratégies d'élagage. Les éléments retenus seront insérés dans l'ensemble \mathcal{Cand}_n . Les stratégies d'élagage adoptées correspondent à :

(i) **L'élagage de tout candidat inclus dans un motif corrélé maximal fréquent**, puisqu'il sera corrélé fréquent d'après la propriété de l'idéal d'ordre des motifs corrélés fréquents (la contrainte de corrélation étant anti-monotone).

(ii) **L'élagage de tout candidat non inclus dans un motif corrélé maximal rare**, puisqu'il ne sera pas corrélé.

(iii) **L'élagage par rapport à la propriété d'idéal d'ordre des motifs minimaux corrélés** : en effet, les motifs minimaux corrélés vérifient la propriété de l'idéal d'ordre. Ainsi, tout candidat minimal corrélé possédant un sous-ensemble non minimal corrélé, sera élagué vu qu'il ne sera pas un motif minimal corrélé.

Notons que tout candidat potentiel inclus dans un motif corrélé maximal rare est forcément corrélé. Toutefois, nous ne pouvons rien confirmer quant à sa rareté. À cet égard, il sera retenu dans l'ensemble \mathcal{Cand}_n et son statut de fréquence sera vérifié grâce à la procédure EXTRACTION_MMCR_MFCR, (cf. ligne 13), dont le pseudo code est donné par l'algorithme 2. Cette procédure permet de déterminer les motifs minimaux corrélés rares à partir des candidats retenus dans l'ensemble \mathcal{Cand}_n . Pour cela, la valeur de *bond* de chaque candidat sera comparée à celles de ses sous-ensembles directs pour déterminer s'il est minimal dans sa classe

d'équivalence ou non. En effet, tout candidat ayant la même valeur de corrélation qu'un de ses sous-ensembles n'est pas minimal de sa classe. Les motifs minimaux corrélés rares identifiés seront ainsi insérés dans l'ensemble \mathcal{MMCR} . Une fois les minimaux repérés, leurs fermetures sont calculées et insérées dans l'ensemble \mathcal{MFCR} . Par ailleurs, dans l'ensemble \mathcal{Cand}_n , seuls les candidats minimaux de leurs classes d'équivalence seront maintenus. Ceci permet d'utiliser \mathcal{Cand}_n dans l'élagage des candidats potentiels de taille $(n + 1)$ (cf. la stratégie d'élagage (iii) de la ligne 11).

Exemple 5 *Considérons la base de transactions donnée par la table 1. L'algorithme RCPR-MINER se déroule de la manière suivante pour $\text{minsupp} = 3$ et $\text{minbond} = 0,20$. Nous avons, initialement, l'ensemble $\mathcal{MCMax} = \{(ACD, 1, \frac{1}{4}), (ABCE, 2, \frac{2}{5})\}$. Étant donné que tous les motifs de cet ensemble sont rares, nous avons donc $\mathcal{MCMaxR} = \{(ACD, 1, \frac{1}{4}), (ABCE, 2, \frac{2}{5})\}$. Ensuite, nous avons $\mathcal{CandP}_1 = \{A, B, C, D, E\}$. Il en dérive, $\mathcal{MMCR}_1 = \{(D, 1, \frac{1}{4})\}$ et $\mathcal{MFCR}_1 = \{(D, 1, \frac{1}{4})\}$. L'ensemble \mathcal{CandP}_2 est ensuite généré : $\mathcal{CandP}_2 = \{AB, AC, AD, AE, BC, BE, BD, CE, CD, DE\}$. Suite à l'application des stratégies d'élagage, nous avons, $\mathcal{MMCR}_2 = \{(AB, 2, \frac{2}{5}), (AE, 2, \frac{2}{5}), (AD, 1, \frac{1}{3}), (CD, 1, \frac{1}{4})\}$. Les motifs fermés associés à ces minimaux, à savoir $(AD, 1, \frac{1}{3})$, $(ACD, 1, \frac{1}{4})$ et $(ABCE, 2, \frac{2}{5})$, sont alors ajoutés à \mathcal{MFCR} . Dans la troisième itération, nous avons $\mathcal{CandP}_3 = \{ABC, ABD, ABE, ACD, ACE, ADE, BCD, BCE, CDE\}$. Aucun de ces candidats n'est minimal rare corrélé. Ainsi, $\mathcal{MMCR}_3 = \{\emptyset\}$. L'ensemble des candidats \mathcal{CandP}_4 est par conséquent vide. Ainsi, les itérations prennent fin donnant ainsi comme résultat les motifs minimaux corrélés rares $\mathcal{MMCR} = \{(D, 1, \frac{1}{4}), (AB, 2, \frac{2}{5}), (AE, 2, \frac{2}{5}), (AD, 1, \frac{1}{3}), (CD, 1, \frac{1}{4})\}$ et leurs fermés $\mathcal{MFCR} = \{(D, 1, \frac{1}{4}), (AD, 1, \frac{1}{3}), (ACD, 1, \frac{1}{4}), (ABCE, 2, \frac{2}{5})\}$.*

3.2 Preuves théoriques

Nous démontrons, dans ce qui suit, les propriétés théoriques de validité et de terminaison de l'algorithme RCPRMINER.

Proposition 1 *L'algorithme RCPRMINER génère tous les motifs minimaux et fermés corrélés rares munis de leurs supports conjonctifs et de leurs valeurs de la mesure bond.*

Preuve. L'algorithme RCPRMINER est un algorithme par niveau permettant d'extraire avec exactitude tous les éléments de la représentation \mathcal{RMCR} . En effet, lors de la première étape, les motifs corrélés maximaux sont identifiés puis ils sont répartis suivant leur statut de fréquence en des motifs corrélés maximaux fréquents et des motifs corrélés maximaux rares. Ces ensembles de motifs seront utilisés pour l'élagage des candidats. Ensuite, les motifs minimaux corrélés rares de l'ensemble \mathcal{MMCR} seront extraits et leurs fermés respectifs seront calculés et insérés dans l'ensemble \mathcal{MFCR} d'une manière itérative.

En effet, lors de chaque itération, un ensemble de candidats de taille n est généré à partir des candidats de taille $n - 1$. Chaque motif candidat doit être inclus dans un motif corrélé maximal rare et ne doit posséder aucun sous-ensemble non minimal corrélé. Ensuite, les supports conjonctifs, disjonctifs, les fermetures conjonctives et les fermetures disjonctives de tous les candidats seront calculés moyennant un balayage du contexte d'extraction. La valeur de la mesure *bond* est ensuite calculée pour tous les candidats retenus. Par la suite, tout candidat possédant un sous-ensemble de même mesure *bond* que lui sera élagué, vu qu'il n'est pas minimal corrélé.

Algorithme 1 : RCPRMINER

Données : Une base de transactions $\mathcal{D} = (\mathcal{T}, \mathcal{I}, \mathcal{R})$, $minbond$, et $minsupp$.
Résultats : La représentation concise exacte $\mathcal{R.MCR} = \mathcal{M.MCR} \cup \mathcal{M.FCR}$.

```

1 Début
2    $\mathcal{R.MCR} := \emptyset$ ;  $Cand_0 := \{\emptyset\}$ ;
3   /* La première étape */
4    $\mathcal{MCMa}x := \text{EXTRACTION\_MCMa}x(\mathcal{D}, minbond)$ ;
5   /* La deuxième étape */
6    $\mathcal{MCMa}xF := \{X \in \mathcal{MCMa}x \mid X.SConj \geq minsupp\}$  /*  $X.SConj$  correspond au
   support conjonctif de  $X$  */;
7    $\mathcal{MCMa}xR := \{X \in \mathcal{MCMa}x \mid X.SConj < minsupp\}$ ;
8    $CandP_1 := \{i \mid i \in \mathcal{I}\}$  /*  $CandP_n$  correspond aux candidats potentiels de taille  $n$  */;
9   tant que ( $CandP_n \neq \emptyset$ ) faire
10    /* Élagage des candidats potentiels */
11     $Cand_n := CandP_n \setminus \{X_n \in CandP_n \mid (\exists Z \in \mathcal{MCMa}xF : X_n \subseteq Z) \text{ ou } (\nexists Z \in$ 
       $\mathcal{MCMa}xR : X_n \subseteq Z) \text{ ou } (\exists Y_{n-1} \subset X_n : Y_{n-1} \notin Cand_{n-1})\}$ ;
12    /* Détermination des motifs minimaux corrélés rares de taille  $n$  et calcul de leurs
      fermetures */
13     $\mathcal{R.MCR} := \mathcal{R.MCR} \cup \text{EXTRACTION\_MMCR\_MFCR}(\mathcal{D}, Cand_n, minsupp)$ ;
14     $n := n + 1$ ;
15     $CandP_n := \text{APRIORI\_GEN}(Cand_{n-1})$ ;
16  retourner  $\mathcal{R.MCR}$ ;
17 Fin

```

L'ensemble des candidats englobe, à ce niveau, tous les motifs minimaux corrélés. Ainsi, chaque candidat rare sera inséré dans l'ensemble $\mathcal{M.MCR}_n$ des motifs minimaux corrélés rares de taille n . Par conséquent, le motif fermé par f_{bond} correspondant au motif minimal corrélé rare en cours, sera calculé. Il résulte, en effet, de l'intersection entre son fermé conjonctif et son fermé disjonctif. Étant donné que les supports conjonctifs, disjonctifs et la mesure $bond$ d'un fermé sont égaux à ceux du motif minimal correspondant, alors nous déduisons que les caractéristiques de chaque fermé par l'opérateur de fermeture f_{bond} sont attribués d'une manière exacte. Ainsi, l'ensemble $\mathcal{M.MCR}_n$ ne contient que les motifs minimaux corrélés rares de taille n et l'ensemble $\mathcal{M.FCR}_n$ ne contient que les fermés corrélés rares de taille n .

L'algorithme marque sa fin d'exécution lorsqu'il n'y a plus de motifs candidats à générer. À la fin de cette étape l'ensemble $\mathcal{M.MCR}$ est composé de tous les motifs qui sont minimaux corrélés rares et leurs fermés respectifs sont inclus dans l'ensemble $\mathcal{M.FCR}$.

Nous concluons que l'algorithme RCPRMINER permet d'extraire avec exactitude tous les éléments des ensembles $\mathcal{M.MCR}$ et $\mathcal{M.FCR}$ munis de leurs supports conjonctifs et de leurs valeurs de la mesure $bond$. Cet algorithme est donc valide et complet.

Proposition 2 *L'algorithme RCPRMINER se termine correctement.*

Preuve. Le nombre des motifs générés par RCPRMINER est fini. En effet, le nombre de motifs candidats pouvant être générés à partir d'un contexte d'extraction ayant n items distincts, est égal au plus à 2^n . De plus, le nombre d'opérations effectuées, afin de traiter chaque candidat est fini. Par conséquent, l'algorithme RCPRMINER se termine correctement.

Algorithme 2 : EXTRACTION_MMCR_MFCR

Données : La base de transactions \mathcal{D} , l'ensemble $Cand_n$ des motifs candidats de taille n , et le seuil minimal de support $minsupp$.

Résultats : L'ensemble \mathcal{MMCR}_n des motifs minimaux corrélés rares de taille n et l'ensemble \mathcal{MFCR} des motifs fermés corrélés rares. L'ensemble $Cand_n$ contenant uniquement les motifs minimaux corrélés.

```

1 Début
2   pour chaque (Transaction  $T$  de  $\mathcal{D}$ ) faire
3     pour chaque ( $X_n \in Cand_n$ ) faire
4        $\omega := X_n \cap T$  /*  $X$  corresponds aux items constituant la transaction  $T$  */;
5       si ( $\omega = \emptyset$ ) alors
6          $X_n.CmpDisj := X_n.CmpDisj \cup T$  /*  $X_n.CmpDisj$  englobe les items
          qui apparaissent dans les transactions ne contenant aucun item de  $X_n$ , par
          conséquent, ces items n'appartiennent donc pas à la fermeture disjonctive du
          candidat  $X_n$ . */;
7        $X_n.SDisj := X_n.SDisj + 1$  /*  $X_n.SDisj$  correspond au support disjonctif de  $X_n$  */;
8       si ( $\omega = X_n$ ) alors
9          $X_n.SConj := X_n.SConj + 1$ ;
10        si  $X_n.f_c = \emptyset$  alors
11           $X_n.f_c := \omega$ ;
12         $X_n.f_c := X_n.f_c \cap \omega$ ;
13   pour chaque ( $X_n \in Cand_n$ ) faire
14      $X_n.bond := \frac{X_n.SConj}{X_n.SDisj}$  /*  $X_n$  est forcément corrélé puisqu'il est inclus dans un motif
      corrélé maximal */;
15     si ( $\exists Y_{n-1} \subset X_n \mid bond(Y_{n-1}) = bond(X_n)$ ) alors
16        $Cand_n := Cand_n \setminus \{X_n\}$  /*  $X_n$  n'est pas un motif minimal corrélé, il est donc
        élagué de l'ensemble  $Cand_n$  et ne sera plus utilisé pour la génération de nouveaux
        candidats */;
17     si ( $X_n.SConj < minsupp$ ) alors
18       /*  $X_n$  est un motif minimal corrélé rare */
19        $\mathcal{MMCR}_n := \mathcal{MMCR}_n \cup (X_n, X_n.SConj, X_n.bond)$ ;
20        $X_n.f_d := T \setminus X_n.CmpDisj$ ;
21        $X_n.f_{bond} := X_n.f_d \cap X_n.f_c$ ;
22        $l := |X_n.f_{bond}|$ ;
23        $\mathcal{MFCR}_l := \mathcal{MFCR}_l \cup (X_n.f_{bond}, X_n.SConj, X_n.bond)$ ;
24        $\mathcal{MFCR} := \mathcal{MFCR} \cup \mathcal{MFCR}_l$ ;
25   retourner ( $\mathcal{MMCR}_n \cup \mathcal{MFCR}$ );
26 Fin

```

Ainsi, nous avons démontré les propriétés théoriques de validité et de terminaison de l'algorithme RCPMINER d'extraction de la représentation \mathcal{RMCR} . Dans la section suivante, nous introduisons l'algorithme ESTMCR d'interrogation de cette représentation.

4 Algorithme ESTMCR d'interrogation de \mathcal{RMCR}

L'interrogation de la représentation permet de déterminer pour un motif donné s'il est corrélé rare. Si c'est le cas, alors les valeurs de son support conjonctif, disjonctif, négatif, ainsi que la valeur de sa mesure *bond*, seront régénérées grâce à la représentation \mathcal{RMCR} . Ceci est réalisé moyennant l'algorithme ESTMCR dont le pseudo-code est donné par l'algorithme 3.

L'algorithme ESTMCR distingue trois différents cas. Le premier se réalise lorsque le motif considéré appartient à la représentation \mathcal{RMCR} . Son support disjonctif et son support négatif seront ainsi aisément dérivés (cf. lignes 3–4). Le deuxième cas se présente lorsque le motif X n'appartient pas à la représentation \mathcal{RMCR} mais il est compris entre deux éléments de cette représentation (cf. ligne 7). Ainsi, le motif fermé associé au motif X correspond au plus petit sur-ensemble, selon l'inclusion ensembliste, appartenant à la représentation \mathcal{RMCR} (cf. ligne 8). Le motif X partage ainsi les mêmes valeurs des différents supports et de *bond* que son fermé (cf. ligne 9–12). Dans le troisième et dernier cas, le motif X n'appartient pas à \mathcal{RMCR} et n'est pas compris entre deux éléments de \mathcal{RMCR} . Ce motif n'est en conséquent pas corrélé rare et l'algorithme retourne un résultat vide (cf. ligne 15). Nous illustrons dans la suite par un exemple l'exécution de l'algorithme ESTMCR.

Exemple 6 Soit la représentation \mathcal{RMCR} donnée par l'exemple 4 (cf. page 5). Considérons le motif ACE . Nous avons $AE \subseteq ACE$ et $ACE \subseteq ABCE$. Ainsi, le motif ACE est corrélé rare. Par ailleurs, sa fermeture est $ABCE$. Par conséquent, $ACE.SConj = ABCE.SConj = 2$, $ACE.SDisj = ABCE.SDisj = 5$, $ACE.SNeg = |T| - ACE.SDisj = 5 - 5 = 0$ et $ACE.bond = ABCE.bond = \frac{2}{5}$. Considérons le motif BC , ce dernier n'appartient pas à \mathcal{RMCR} et il n'est pas compris entre deux éléments de la représentation. Ainsi, l'algorithme ESTMCR retourne un résultat vide pour indiquer que le motif BC n'est pas un motif corrélé rare.

5 Algorithme REGENERATIONMCR de régénération de \mathcal{MCR}

La régénération de l'ensemble \mathcal{MCR} à partir de \mathcal{RMCR} s'effectue grâce à l'algorithme REGENERATIONMCR dont le pseudo-code est donné par l'algorithme 4. Cet algorithme fournit l'ensemble \mathcal{MCR} des motifs corrélés rares munis de leurs supports conjonctifs et de leurs valeurs de la mesure *bond*. L'exemple suivant illustre l'exécution de cet algorithme.

La tâche de régénération s'effectue à travers l'algorithme REGENERATIONMCR de la manière suivante. D'abord, tous les éléments de la représentation \mathcal{RMCR} seront insérés dans l'ensemble \mathcal{MCR} (cf. ligne 4) initialement vide. Par la suite, l'algorithme parcourt l'ensemble \mathcal{MMCR} des motifs minimaux et affecte à chaque motif minimal M son fermé F (cf. ligne 6). Puis l'ensemble de motifs compris entre le minimal M et son fermé F est généré (cf. ligne 7). Chaque élément de cet ensemble est un motif corrélé rare et partage le même support conjonctif et la même valeur de *bond* que son fermé F et sera inséré dans l'ensemble \mathcal{MCR} (cf. ligne 10). Lorsque tous les motifs générés sont insérés dans l'ensemble \mathcal{MCR} , alors l'algorithme retourne l'ensemble total des motifs corrélés rares \mathcal{MCR} (cf. ligne 11).

Exemple 7 Considérons la représentation concise exacte donnée par l'exemple 4 (cf. page 5). D'abord, l'ensemble \mathcal{MCR} est initialisé par l'algorithme REGENERATIONMCR à l'ensemble vide. Tous les éléments de \mathcal{RMCR} seront ensuite insérés dans l'ensemble \mathcal{MCR} . Ainsi, \mathcal{MCR}

Algorithme 3 : ESTMCR

Données : La représentation $\mathcal{RMCR} = \mathcal{MMCR} \cup \mathcal{MFRCR}$, un motif X , et le nombre de transactions de la base, c.-à.-d., $|\mathcal{T}|$.

Résultats : Le support conjonctif, disjonctif, négatif et la valeur de la mesure *bond* si le motif X est corrélé rare. Sinon, un résultat vide est retourné.

1 **Début**

2 **si** ($X \in \mathcal{RMCR}$) **alors**

3 $X.SDisj = \frac{X.SConj}{X.bond}$ /* $X.SDisj$ correspond au support disjonctif de X */;

4 $X.SNeg = |\mathcal{T}| - X.SDisj$ /* $X.SNeg$ correspond au support négatif de X */;

5 **retourner** $\{X, X.SConj, X.SDisj, X.SNeg, X.bond\}$;

6 **sinon**

7 **si** ($\exists Y, Z \in \mathcal{RMCR} \mid Y \subset X \text{ et } X \subset Z$) **alors**

8 $F := \min_{\subseteq} \{X_1 \in \mathcal{RMCR} \mid X \subset X_1\}$ /* F dénote la fermeture de X , repérée étant le plus petit motif par inclusion ensembliste de la représentation englobant X */;

9 $X.SConj = F.SConj$;

10 $X.bond = F.bond$;

11 $X.SDisj = \frac{X.SConj}{X.bond}$;

12 $X.SNeg = |\mathcal{T}| - X.SDisj$;

13 **retourner** $\{X, X.SConj, X.SDisj, X.SNeg, X.bond\}$;

14 **sinon**

15 **retourner** \emptyset ;

16 **Fin**

$= \{(D, 1, \frac{1}{1}), (AB, 2, \frac{2}{5}), (AD, 1, \frac{1}{3}), (AE, 2, \frac{2}{5}), (CD, 1, \frac{1}{4}), (ACD, 1, \frac{1}{4}), (ABCE, 2, \frac{2}{5})\}$. Par la suite, nous générons les motifs ABE et ABC compris entre le minimal $(AB, 2, \frac{2}{5})$ et son fermé $(ABCE, 2, \frac{2}{5})$ et le motif ACE compris entre le minimal $(AE, 2, \frac{2}{5})$ et son fermé $(ABCE, 2, \frac{2}{5})$. Les motifs ABE , ABC et ACE générés seront alors insérés dans l'ensemble \mathcal{MCR} . Ce dernier englobe, ainsi, tous les motifs corrélés rares. $\mathcal{MCR} = \{(D, 1, \frac{1}{1}), (AB, 2, \frac{2}{5}), (AD, 1, \frac{1}{3}), (AE, 2, \frac{2}{5}), (CD, 1, \frac{1}{4}), (ABC, 2, \frac{2}{5}), (ABE, 2, \frac{2}{5}), (ACD, 1, \frac{1}{4}), (ACE, 2, \frac{2}{5}), (ABCE, 2, \frac{2}{5})\}$.

6 Évaluation expérimentale de la représentation \mathcal{RMCR}

Notre objectif principal, dans cette section, est de prouver expérimentalement le taux de compacité de la représentation \mathcal{RMCR} . Les différentes expérimentations réalisées ont été menées sur une machine munie d'un processeur Intel Dual Core E5400, ayant une fréquence de 2,7GHz avec 4Go de mémoire vive, tournant sur une plateforme Linux Ubuntu 10.04. Les expérimentations ont été réalisées sur différentes bases de test benchmark denses et éparses ⁽⁴⁾.

Les résultats expérimentaux les plus représentatifs sont donnés par la figure 3. Nous constatons, d'après ces résultats, que les taux de réduction obtenus pour la représentation proposée et pour différents seuils *minsupp* et *minbond* sont intéressants. Par ailleurs, la représentation

4. Disponibles à l'adresse suivante : <http://fimi.cs.helsinki.fi/data>.

Algorithme 4 : REGENERATIONMCR

Données : La représentation concise exacte $\mathcal{RMCR} = \mathcal{MMCR} \cup \mathcal{MFCR}$.

Résultats : L'ensemble \mathcal{MCR} des motifs corrélés rares munis de leurs valeurs du support conjonctif et de leurs valeurs de la mesure *bond*.

```

1 Début
2    $\mathcal{MCR} := \emptyset$ ;
3   pour chaque ( $X \in \mathcal{RMCR}$ ) faire
4      $\mathcal{MCR} := \mathcal{MCR} \cup \{X, X.SConj, X.bond\}$  ;
5   pour chaque ( $M \in \mathcal{MMCR}$ ) faire
6      $F := \min_{\subseteq} \{M_1 \in \mathcal{MFCR} \mid M \subset M_1\}$  /*  $F$  dénote la fermeture du motif minimal
       corrélé rare  $M$ , repérée étant le plus petit motif par inclusion ensembliste de la
       représentation englobant  $M$  */ ;
7     pour chaque ( $X \mid M \subset X$  et  $X \subset F$ ) faire
8        $X.SConj = F.SConj$ ;
9        $X.bond = F.bond$ ;
10     $\mathcal{MCR} := \mathcal{MCR} \cup \{X, X.SConj, X.bond\}$  ;
11  retourner  $\mathcal{MCR}$ ;
12 Fin
    
```

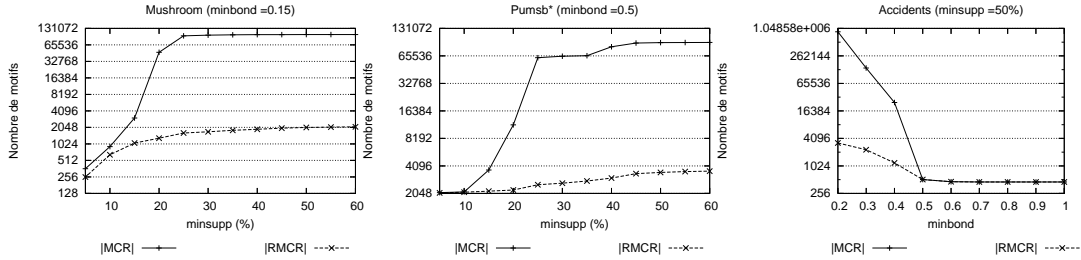


FIG. 3 – Variation des cardinalités de la représentation \mathcal{RMCR} versus celles de l'ensemble \mathcal{MCR} en fonction de \minsupp et de \minbond .

\mathcal{RMCR} est prouvée être une couverture parfaite de l'ensemble \mathcal{MCR} dans le sens que sa taille ne dépasse jamais celle de ce dernier. Considérons par exemple, la base MUSHROOM pour $\minsupp = 35\%$ et $\minbond = 0,15$: $|\mathcal{RMCR}| = 1\,810$, et $|\mathcal{MCR}| = 100\,156$. Le taux de compacité dans ce cas est de 98%. Ces résultats sont obtenus grâce à la propriété de non-injectivité de l'opérateur de fermeture f_{bond} . En effet, cet opérateur permet de regrouper les motifs ayant les mêmes propriétés dans une même classe d'équivalence. Ceci permet ainsi d'éviter la redondance des éléments maintenus. Nous avons, par exemple, pour la base MUSHROOM : $|\mathcal{MMCR}| = 1\,412$ et $|\mathcal{MFCR}| = 652$. Puisque la représentation \mathcal{RMCR} correspond à l'union sans redondance des ensembles \mathcal{MMCR} et \mathcal{MFCR} , nous avons toujours $|\mathcal{RMCR}| \leq |\mathcal{MMCR}| + |\mathcal{MFCR}|$.

Dans ce qui suit, nous proposons une application de \mathcal{RMCR} dans le cadre de la détection d'intrusions dans les réseaux informatiques.

7 Application de la représentation \mathcal{RMCR} dans la détection d'intrusions

Nous présentons dans cette section, l'application de la représentation \mathcal{RMCR} dans un processus de classification basé sur les règles d'association corrélées rares. En effet, les ensembles de motifs \mathcal{MMCR} et \mathcal{MFCR} , composant la représentation \mathcal{RMCR} , seront incorporés dans la dérivation des règles d'association génériques corrélées rares de la forme $Gen \Rightarrow Fermé \setminus Gen$, avec $Gen \in \mathcal{MMCR}$ et $Fermé \in \mathcal{MFCR}$ ⁽⁵⁾.

Ensuite, à partir des règles génériques obtenues, nous extrayons les règles de classification. En effet, les règles génériques obtenues seront filtrées afin de ne garder que les règles génériques ayant le libellé de la classe d'attaque dans la partie conclusion. Ces règles seront alors communiquées au classifieur que nous avons conçu. Ce dernier permet d'élaborer la tâche de classification et retourne le taux de détection pour chaque classe d'attaque. Nous présentons dans la suite l'évaluation expérimentale de la classification basée sur les règles corrélées rares pour la base de données KDD 99 ⁽⁶⁾.

7.1 Description de la base KDD 99

Chaque objet de la base KDD 99 représente une connexion du flot de données. Une connexion est ainsi étiquetée comme *Normale* ou *Attaque*. La base KDD 99 décrit 38 catégories d'attaques réparties en quatre classes d'attaques, à savoir DOS, PROBE, R2L et U2R, et une classe NORMALE. Cette base contient 4 940 190 objets dans la base d'apprentissage et chaque objet est caractérisé par 41 attributs. Nous considérons, dans ce travail, 10% de l'ensemble d'apprentissage dans la phase de construction du classifieur, contenant ainsi 494 019 objets. L'ensemble d'apprentissage contient 79,20% (respectivement, 0,83%, 0,22% et 0,10%) d'attaques DOS (respectivement, PROBE, R2L et U2R), et le reste, *c.-à.-d.* 19,65%, concerne des connexions étiquetées *Normale*.

7.2 Discussion des résultats obtenus

Les résultats expérimentaux obtenus sont donnés par la table 2, avec "RAs" et "TD" les abréviations respectives de "Règles d'Association" et "Taux de Détection", et *minconf* dénote le seuil minimal de la mesure *confiance* (Agrawal et Srikant, 1994). Nous entendons aussi par "Phase de construction" l'étape de l'extraction de la représentation \mathcal{RMCR} tandis que par "Phase de classification", nous entendons l'étape de dérivation des règles de classification à partir de la représentation \mathcal{RMCR} et leur application dans la détection d'intrusions.

Nous constatons que les taux de détection les plus intéressants sont achevés pour les classes d'attaques NORMALE et DOS. En effet, ceci est expliqué par la taille élevée en nombre de connexions de ces deux classes d'attaques. Ceci confirme que notre approche proposée dans ce travail présente de meilleures performances pour des bases volumineuses. Nous remarquons aussi que ce taux de détection varie d'une classe d'attaque à une autre. Par exemple, pour la classe U2R, ce taux est relativement faible par rapport aux autres classes d'attaques.

5. Par "générique", nous entendons que ces règles sont à prémisse minimale et à conclusion maximale, selon la relation d'inclusion ensembliste.

6. La base KDD 99 est disponible à l'adresse suivante : <http://kdd.ics.uci.edu/databases/kddcup99/kddcup99.html>.

Fouille d'une représentation concise des motifs corrélés rares

Nous concluons aussi, d'après les résultats de la table 2, que les coûts de calcul varient d'une classe d'attaque à une autre. Toutefois, pour les différentes classes d'attaques considérées, la phase de construction est plus coûteuse en temps d'exécution que la phase de classification. Ceci est justifié par le fait que l'étape de construction englobe l'extraction de la représentation concise \mathcal{RMCR} , or cette opération est NP-difficile (Boley et Gärtner, 2009) étant donnée la complexité liée à la localisation des deux bordures associées aux contraintes de corrélation et de rareté.

Classe de l'attaque	<i>minsupp</i> (%)	<i>minbond</i>	<i>minconf</i>	# RAs génériques exactes	# RAs génériques approximatives	# RAs génériques de classification	TD (%)	Temps CPU (en secondes)	
								Phase de construction	Phase de classification
DOS	80	0,95	0,90	4	31	17	98,68	120	1
PROBE	60	0,70	0,90	232	561	15	70,69	55	1
R2L	80	0,90	0,70	2	368	1	81,52	1 729	1
U2R	60	0,75	0,75	106	3	5	38,46	32	1
NORMALE	85	0,95	0,95	0	10	3	100,00	393	15

TAB. 2 – Évaluation des règles d'association corrélées rares pour la base KDD 99.

La table 3 compare les résultats obtenus par notre approche, basée sur les règles d'association corrélées rares, à ceux offerts par les approches basées respectivement sur les arbres de décisions et les réseaux bayésiens (Ben Amor et al., 2004). Il est à noter que le choix de ces approches pour ce comparer avec est argumenté par le fait que celle utilisant les arbres de décisions est aussi basée sur les règles d'association. Par ailleurs, l'apprentissage est supervisé dans les différentes approches comparées. Les résultats obtenus prouvent que notre approche offre dans différentes situations de meilleures performances que les autres approches. En effet, elle est la meilleure pour les classes d'attaques DOS, R2L et U2R. Bien que aussi meilleurs pour la classe NORMALE, les résultats obtenus sont très proches de ceux obtenus avec les arbres de décision. Les réseaux bayésiens présentent de meilleurs taux de détection uniquement pour la classe PROBE. Ainsi, l'application des règles corrélées rares offre une solution intéressante dans le contexte de la détection d'intrusions.

Classe d'attaque	RAs corrélées rares	Arbres de décision	Réseaux bayésiens
DOS	98,68	97,24	96,65
PROBE	70,69	77,92	88,33
R2L	81,52	0,52	8,66
U2R	38,46	13,60	11,84
NORMALE	100,00	99,50	97,68

TAB. 3 – Comparaison des taux de détection obtenus pour les règles corrélées rares versus les approches de l'état de l'art.

8 Conclusion et perspectives

Dans ce papier, nous avons proposé l'algorithme RCPRMINER d'extraction de la représentation concise exacte \mathcal{RMCR} de l'ensemble \mathcal{MCR} des motifs corrélés rares. Nous avons introduit également l'algorithme ESTMCR d'interrogation de cette représentation ainsi que l'algorithme REGENERATIONMCR de dérivation de l'ensemble \mathcal{MCR} à partir de \mathcal{RMCR} .

Nous avons démontré expérimentalement le taux de réduction intéressant offert par cette représentation. L'efficacité de la classification, basée sur les règles d'association corrélées rares, a été aussi prouvée dans le cadre de la détection d'intrusions.

Les perspectives de travaux futurs concernent : (i) La comparaison détaillée des performances d'un algorithme d'extraction de \mathcal{MCR} , directement à partir d'une base de transactions, à celles de RCPRMINE suivi par REGENERATIONMCR pour dériver l'ensemble total des motifs corrélés rares à partir de \mathcal{RMCR} . Ceci permettra de cerner aussi les situations où le recours à la représentation \mathcal{RMCR} est aussi nécessaire non seulement pour réduire la taille des connaissances extraites mais aussi pour rendre possible la fouille des motifs corrélés rares. (ii) L'extraction, à partir de \mathcal{RMCR} , de formes généralisées de règles d'association présentant des conjonctions, des disjonctions, et des négations d'items en prémisse ou en conclusion ainsi que leur application dans des contextes réels. (iii) L'extension de l'approche proposée pour d'autres mesures de corrélation (Kim et al., 2011; Omiecinski, 2003; Segond et Borgelt, 2011; Surana et al., 2010; Xiong et al., 2006) en se basant sur l'étude de leurs propriétés respectives.

Références

- Agrawal, R. et R. Srikant (1994). Fast algorithms for mining association rules in large databases. In *Proceedings of the 20th International Conference on Very Large Data Bases (VLDB 1994)*, Santiago, Chile, pp. 487–499.
- Ben Amor, N., S. Benferhat, et Z. Elouedi (2004). Naive bayes vs decision trees in intrusion detection systems. In *Proceedings of the ACM Symposium on Applied Computing (SAC 2004)*, Nicosia, Cyprus, 2004, pp. 420–424.
- Ben Younes, N., T. Hamrouni, et S. Ben Yahia (2012). À la recherche des motifs corrélés : proposition d'une nouvelle représentation concise exacte associée à la mesure *bond*. À paraître dans la revue *Technique et Science Informatiques (TSI)*.
- Boley, M. et T. Gärtner (2009). On the complexity of constraint-based theory extraction. In *Proceedings of the 12th International Conference Discovery Science (DS 2009)*, LNCS, volume 5808, Springer-Verlag, Porto, Portugal, pp. 92–106.
- Bouasker, S., T. Hamrouni, et S. Ben Yahia (2012). Motifs corrélés rares : caractérisation et nouvelles représentations concises exactes. À paraître dans le numéro spécial de la RNTI : “Mesurer et évaluer la qualité des données et des connaissances”.
- Brin, S., R. Motwani, et C. Silverstein (1997). Beyond market baskets : generalizing association rules to correlations. In *Proceedings of the ACM-SIGMOD International Conference on Management of Data (SIGMOD 1997)*, Washington D. C., USA, pp. 265–276.
- Cohen, E., M. Datar, S. Fujiwara, A. Gionis, P. Indyk, R. Motwani, J. D. Ullman, et C. Yang (2000). Finding interesting associations without support pruning. In *Proceedings of the 16th International Conference on Data Engineering (ICDE 2000)*, IEEE Computer Society Press, San Diego, California, USA, pp. 489–499.
- Grahne, G., L. V. S. Lakshmanan, et X. Wang (2000). Efficient mining of constrained correlated sets. In *Proceedings of the 16th International Conference on Data Engineering (ICDE 2000)*, IEEE Computer Society Press, San Diego, California, USA, pp. 512–521.
- Kim, S., M. Barsky, et J. Han (2011). Efficient mining of top correlated patterns based on null-invariant measures. In *Proceedings of the 15th European Conference on Machine Learning and Knowledge Discovery in Databases (ECML/PKDD 2011)*, LNCS, volume 6912, Springer, Athens, Greece, pp. 177–192.

Ma, S. et J. L. Hellerstein (2001). Mining mutually dependent patterns. In *Proceedings of the 1st International Conference on Data Mining (ICDM 2001)*, IEEE Computer Society Press, San Jose, California, USA, pp. 409–406.

Omiecinski, E. (2003). Alternative interest measures for mining associations in databases. *IEEE Transactions on Knowledge and Data Engineering* 15(1), 57–69.

Sandler, I. et A. Thomo (2010). Mining frequent highly-correlated item-pairs at very low support levels. In *Proceedings of the Workshop on High Performance Analytics - Algorithms, Implementations, and Applications (PHPA 2010) in conjunction with the 10th SIAM International Conference on Data Mining (SDM 2010)*, Columbus, Ohio, USA.

Segond, M. et C. Borgelt (2011). Item set mining based on cover similarity. In *Proceedings of the 15th Pacific-Asia Conference on Knowledge Discovery and Data Mining (PAKDD 2011)*, LNCS, volume 6635, Springer-Verlag, Shenzhen, China, pp. 493–505.

Surana, A., R. U. Kiran, et P. K. Reddy (2010). Selecting a right interestingness measure for rare association rules. In *Proceedings of the 16th International Conference on Management of Data (COMAD 2010)*, Nagpur, India, pp. 115–124.

Xiong, H., P. N. Tan, et V. Kumar (2006). Hyperclique pattern discovery. *Data Mining and Knowledge Discovery*. 13(2), 219–242.

Summary

Correlated rare pattern mining is an interesting issue in Data mining. In this respect, the set of correlated rare patterns w.r.t. to the *bond* correlation measure was studied in a recent work, in which the \mathcal{RCPR} concise exact representation of the set of correlated rare patterns was proposed. However, none algorithm was proposed in order to mine this representation and none experiment was carried out to evaluate it. In this paper, we introduce the new RCPRMINE algorithm allowing an efficient extraction of \mathcal{RCPR} . We also present the ISRCP algorithm allowing the query of the \mathcal{RCPR} representation in addition to the RCPREGENERATION algorithm allowing the regeneration of the whole set \mathcal{RCP} of rare correlated patterns starting from this representation. The carried out experiments highlight interesting compactness rates offered by \mathcal{RCPR} . The effectiveness of the proposed classification method, based on generic rare correlated association rules derived from \mathcal{RCPR} , has also been proved in the context of intrusion detection.